

## Internet traffic volumes are not Gaussian - they are log-normal: an 18-year longitudinal study with implications for modelling and prediction

Article (Accepted Version)

Alasmar, Mohammed, Clegg, Richard, Zakhleniuk, Nickolay and Parisi, George (2021) Internet traffic volumes are not Gaussian - they are log-normal: an 18-year longitudinal study with implications for modelling and prediction. IEEE/ACM Transactions on Networking, 29 (3). pp. 1266-1279. ISSN 1063-6692

This version is available from Sussex Research Online: <http://sro.sussex.ac.uk/id/eprint/96912/>

This document is made available in accordance with publisher policies and may differ from the published version or from the version of record. If you wish to cite this item you are advised to consult the publisher's version. Please see the URL above for details on accessing the published version.

### **Copyright and reuse:**

Sussex Research Online is a digital repository of the research output of the University.

Copyright and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable, the material made available in SRO has been checked for eligibility before being made available.

Copies of full text items generally can be reproduced, displayed or performed and given to third parties in any format or medium for personal research or study, educational, or not-for-profit purposes without prior permission or charge, provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

# Internet Traffic Volumes Are Not Gaussian - They Are Log-Normal: An 18-Year Longitudinal Study With Implications for Modelling and Prediction

Mohammed Alasmar  
Department of Informatics  
University of Sussex  
Brighton, UK  
m.alasmar@sussex.ac.uk

Richard Clegg  
School of Computer Science  
Queen Mary University of London  
London, UK  
r.clegg@qmul.ac.uk

Nickolay Zakhleniuk  
School of Computer Science  
University of Essex  
Colchester, UK  
naz@essex.ac.uk

George Parisi  
Department of Informatics  
University of Sussex  
Brighton, UK  
g.parisi@sussex.ac.uk

**Abstract**—Getting good statistical models of traffic on network links is a well-known, often-studied problem. A lot of attention has been given to correlation patterns and flow duration. The distribution of the amount of traffic per unit time is an equally important but less studied problem. We study a large number of traffic traces from many different networks including academic, commercial and residential networks using state-of-the-art statistical techniques. We show that traffic obeys the log-normal distribution which is a better fit than the Gaussian distribution commonly claimed in the literature. We also investigate an alternative heavy-tailed distribution (the Weibull) and show that its performance is better than Gaussian but worse than log-normal. We examine anomalous traces which exhibit a poor fit for all distributions tried and show that this is often due to traffic outages or links that hit maximum capacity. We demonstrate that the data we look at is stationary if we consider samples of 15-minute long or even 1-hour long. This gives confidence that we can use the distributions for estimation and modelling purposes.

We demonstrate the utility of our findings in two contexts: predicting that the proportion of time traffic will exceed a given level (for service level agreement or link capacity estimation) and predicting 95th percentile pricing. We also show that the log-normal distribution is a better predictor than Gaussian or Weibull distributions in both contexts.

**Index Terms**—Traffic modelling, network planning, bandwidth provisioning, traffic billing

## I. INTRODUCTION

Internet traffic characterisation is an important problem for network researchers and vendors. The subject has a long history. Early works [1], [2] discovered that the correlation structure of traffic exhibits self-similarity and that the durations of individual flows of packets show heavy-tails [3]. These works were later challenged and refined (see Section VII for a summary). By comparison, the distribution of the amount of traffic present on a link in a given time period has seen comparatively less research interest. This is surprising as correct traffic statistics can be extremely useful in network planning. In this paper we use a rigorous statistical approach to fitting a statistical distribution to the amount of traffic within a given time period. Formally, we choose some timescale  $T$  and let  $X_i$  be the amount of traffic seen in the time period  $[iT, (i+1)T)$ . We investigate the distribution of the random

variable  $X$  over a wide range of values of  $T$ . We show that the distribution of the variable has considerable implications for network planning; for assessing how often a link is over capacity and in particular for service level agreements (SLAs), and for traffic pricing, particularly using the 95th percentile scheme [4].

Previous authors have claimed that  $X$  has a normal (or Gaussian) distribution [5]–[7]. Others claim  $X$  is Gaussian plus a tail associated with bursts [8], [9]. All these studies are based on straightforward goodness-of-fit tests (e.g. Quantile-Quantile (Q-Q) plots) and relevant correlation tests that are used to assess how well captured traffic traces are fitted to Gaussian or heavy-tailed distributions. As discussed in [10], these statistical approaches can produce a substantially inaccurate assessment about whether samples follow a Gaussian/heavy-tailed or not. This is because the difference in these distributions lies in the behaviour of the tail where there can be relatively few samples, therefore large amounts of data and careful statistical handling are both important to determine the correct distribution [10].

In this paper, we use a well-established statistical methodology [10] to show that a log-normal<sup>1</sup> distribution is a better fit than Gaussian or Weibull<sup>2</sup> for the vast majority of traces. This holds over a wide range of timescales  $T$  (from 5 ms to 5 sec) [11]. The question becomes whether the log-normal model remains applicable for all traffic volumes that have the same time periods and aggregation timescales used for testing this model. Stationarity tests can answer this question as stationarity from a probabilistic point of view means that the distribution remains unchanged for whatever shift of time window [12]–[14]. The vast majority of modelling techniques developed for volume-based traffic profiling imply the assumption of statistical stationarity which is often taken

<sup>1</sup>A variable  $X$  has a log-normal distribution if its logarithm is normally distributed  $\ln(X) \sim N(\mu, \sigma^2)$  where  $\mu \in \mathbb{R}$  is the mean and  $\sigma > 0$  is the standard deviation of the distribution (see Table 1 in [10]).

<sup>2</sup>A variable  $X$  has a Weibull distribution with parameters  $k > 0$  (known as shape) and  $\lambda > 0$  (known as scale) if its probability density function follows  $f(x) = \frac{k}{\lambda} \left(\frac{x}{\lambda}\right)^{k-1} \exp(-(x/\lambda)^k)$  when  $x \geq 0$  and is 0 otherwise.

for granted without being explicitly validated [5], [6], [8]. In contrast, we extensively test all studied time series for stationarity using state of the art techniques and examine their trend and seasonality components. The majority of the 15-minute and 1-hour long traces in the dataset are stationary when aggregated at timescales of 500 ms to 5 sec.

This paper is the most comprehensive investigation of this phenomenon the authors know about. We study a large number of publicly available traces from a diverse set of locations (including commercial, academic and residential networks) with different link speeds and spanning the last 18 years. There is a small number of anomalous traces in our datasets where the distribution deviates from log-normal and we find that this occurs when a link spends considerable time either having an outage or completely at maximum capacity. These anomalous traces can be presented using a bimodal distribution.

We show how often a link following a given distribution will be over a given capacity and show that our approach improves greatly on results which assume that traffic follows a Gaussian distribution. We further show that if an ISP wishes to estimate future transit bills that use the 95th percentile billing scheme, then the log-normal is a better model than the Gaussian distribution.

The structure of the paper is as follows. In Section II we describe the datasets used. In Section III we describe our best-practice procedure for fitting traffic and demonstrate that log-normal is the best fit distribution among all studied distributions for our traces under a variety of circumstances. We examine those few traces that do not follow this distribution and find it occurs when a link spends considerable time either having an outage or being completely at maximum capacity. In Section IV we test all studied time series for stationarity at different timescales. In Section V we demonstrate that the log-normal distribution is the most suitable for estimating how often a link is over capacity. In Section VI we show that the log-normal distribution provides good estimates when looking at 95th percentile pricing. In Section VII we give related work. Finally, Section VIII gives our conclusions.

## II. NETWORK TRAFFIC TRACES

A key contribution of our work stems from the spatial and temporal diversity of the studied traces. The studied dataset spans a period of 18 years and comprises 232 traces<sup>3</sup>.

**CAIDA traces.** We have used 27 CAIDA traces captured at an Internet data collection monitor which is located at an Equinix data centre in Chicago [15]. The data centre is connected to a backbone link of a Tier 1 ISP. The monitor records an hour-long trace four times a year, usually from 13:00 to 14:00 UTC. The selection of a 15-minute trace from the respective 1-hour trace is done by first splitting the trace into four 15-minute subtraces and, subsequently, selecting one of these at random. Each trace contains billions of IPv4 packets, the headers of which are

anonymised. The average captured data rate is 2.5 Gbps. At the time of capturing, the monitored link had a capacity of 10 Gbps. Traces were captured between 2013 and 2016.

**MAWI traces.** The MAWI archive [16] consists of a collection of Internet traffic traces, captured within the WIDE backbone network that connects Japanese universities and research institutions to the Internet. Each trace consists of IP level traffic observed daily from 14:00 to 14:15 at a vantage point within WIDE. Traces include anonymised IP and MAC headers, along with an *ntpd* timestamp [16]. We have looked at 110 traces (each one being 15-minute long). Traces were captured between 2014 and 2020. On average, each trace consists of 70 million packets; the average captured data rate is 422 Mbps. The monitored link had a capacity of 1 Gbps. For the stationarity tests presented in Section IV we used a 24-hour long MAWI trace<sup>4</sup>. This trace was captured on 09/05/2018 at samplepoint-G which monitors a 10 Gbps link to DIX-IE<sup>5</sup>.

**Twente University traces.** We used 40 traffic traces captured at five different locations (8 traces from each location). Traces are diverse in terms of the link rates, types of users and capture time [17]. Each trace is 15-minute long. The first location is a residential network with a 300 Mbps link, which connects 2000 students (each one having a 100 Mbps access link); traces were captured in July 2002. The second location is a research institute network with a 1 Gbps link which connects 200 researchers (each one having a 100 Mbps access link); traces were captured between May and August 2003. The third location is at a large college with a 1 Gbps link which connects 1000 employees (each one having a 100 Mbps access link); traces were captured between February and July 2004. The fourth location is an ADSL access network with a 1 Gbps ADSL link used by hundreds of users (each one having a 256 Kbps to 8 Mbps access link); traces were captured between February and July 2004. The fifth location is an educational organisation with a 100 Mbps link connecting 135 students and employees (each one having a 100 Mbps access link); traces were captured between May and June 2007.

**Waikato University VIII traces.** The Waikato dataset consists of traffic traces captured by the WAND group at the University of Waikato, New Zealand [18]. The capture point is at the link interconnecting the University with the Internet. All of the traces were captured using software that was specifically developed for the Waikato capture point and a DAG 3 series hardware capture card. All IP addresses within the traces are anonymised. In our study, we have used 30 traces captured between April 2011 and November 2011.

**Auckland University IX traces.** The Auckland dataset consists of traffic traces captured by the WAND group at the University of Waikato [19]. The traces were collected at the University of Auckland, New Zealand. The capture point is at the link interconnecting the University with the Internet. All IP addresses within the traces are anonymised. In our study, we have used 25 traces captured in 2009.

<sup>3</sup>All Matlab and Python scripts used for the data analysis presented in this paper can be found on <https://github.com/mohammedalasar/InternetTraces>

<sup>4</sup><https://mawi.wide.ad.jp/mawi/ditl/ditl2018-G/>

<sup>5</sup><http://two.wide.ad.jp/>

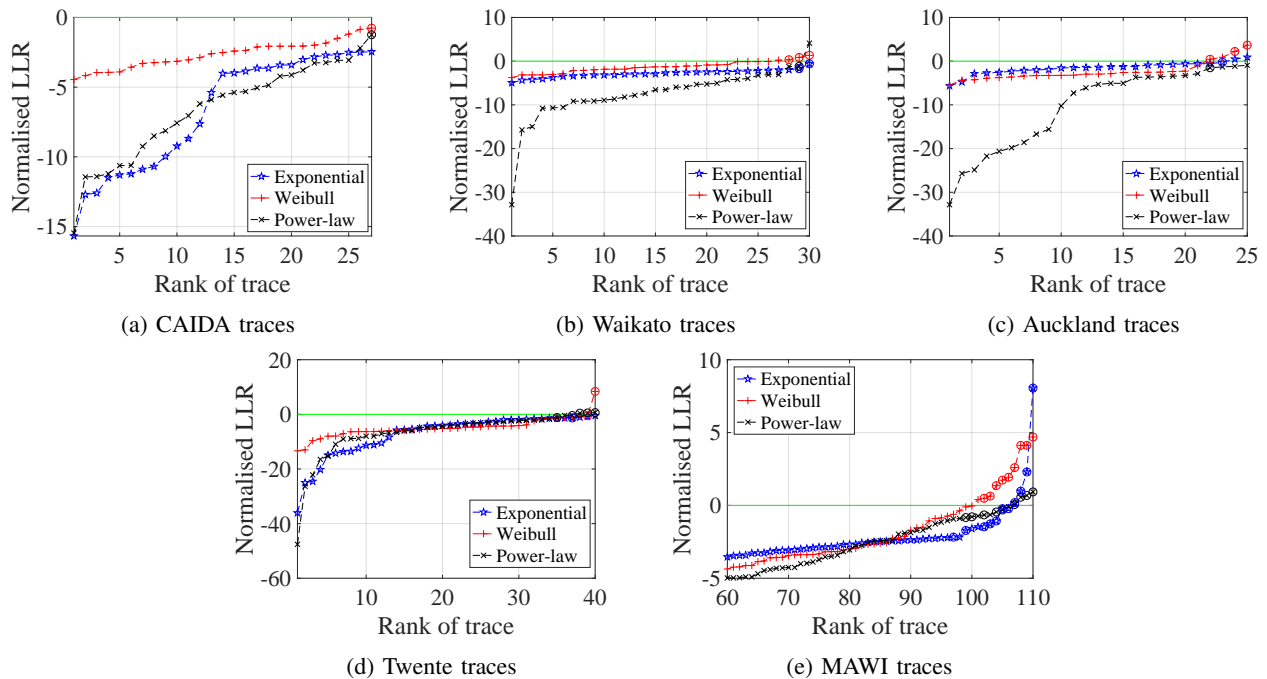


Fig. 1: Normalised Log-Likelihood Ratio ( $\mathfrak{R}$ ) test results for all studied traces and candidate distributions when running *fit.distCompare*(alternative, log-normal). Aggregation timescale  $T$  is 100 msec. Circled points in the plots are the ones with  $p$ -value greater than 0.1, i.e. test is inconclusive with respect to assessing which of the reference and candidate distributions is a better fit to the traffic data.

### III. FITTING A STATISTICAL DISTRIBUTION TO INTERNET TRAFFIC DATA

In this section we present an extensive statistical analysis applied to the datasets described in the previous section. The aim is to find a good model for all studied traces for different aggregation time values. In contrast to the existing research (see Section VII), we are basing our analysis on the framework proposed by Clauset et al. [10], a comprehensive statistical framework developed specifically for testing power-law (or another reference distribution) behaviour in empirical data<sup>6</sup>. The framework combines maximum-likelihood fitting methods with goodness-of-fit tests based on the Kolmogorov–Smirnov statistic and likelihood ratios. The method reliably tests whether the reference distribution is the best model for a specific dataset, or, if not, whether an alternative statistical distribution is. The framework performs the tests described above as follows: (1) the parameters of the reference model are estimated for a given dataset; (2) the goodness-of-fit between the data and the reference distribution is calculated, under the hypothesis that the reference distribution is a good fit to the provided traffic samples. If the resulting  $p$ -value is greater than 0.1 the hypothesis is accepted (i.e. the reference distribution is a plausible fit to the given data), otherwise the hypothesis is rejected; (3) alternative distributions are tested

against the the reference distribution as a better fit to the data by employing a likelihood ratio test.

In this paper, we use the log-normal as the reference distribution and compare it to the exponential, Weibull and power-law distributions (see Equation 1). This is in contrast to our initial results presented in [11], where following Clauset method we used the power-law as the reference distribution. This is because (1) the power-law distribution failed to fit the vast majority of the studied datasets and (2), as shown in [11], the log-normal distribution was a good fit for most of the traces. Indeed, Step 2 of the Clauset method showed that for the vast majority of the examined traces, the respective hypothesis was accepted; i.e. the log-normal distribution was a good fit. This is an encouraging result; here, we build on it by comparing the log-normal to three alternative distributions, exponential, Weibull and power-law, using Step 3 of the Clauset method, i.e. by performing the log-likelihood ratio (LLR) test, as follows:

$$\mathfrak{R}, p = \text{fit.distCompare}(\text{alternative}, \text{lognormal}) \quad (1)$$

where  $\mathfrak{R}$  is the normalised LLR<sup>7</sup> between the alternative and log-normal distributions,  $p$  is the significance value for this test.  $\mathfrak{R}$  is positive if the alternative distribution is a better fit

<sup>7</sup> $\mathfrak{R}$  is calculated as  $\mathcal{R}/(\sigma\sqrt{n})$ , where  $\mathcal{R}$  is the log-likelihood ratio [10]. Note that if we run *fit.distCompare*(Y,X) and *fit.distCompare*(Z,X), then we would not be able to tell which of the Y and Z distributions is a better fit to the used data, even if both were better fits compared to X. This is because the normalised LLR value is measured by normalising LLR by  $\sigma$  (the estimated standard deviation on LLR), which is a nonlinear operator.

<sup>6</sup>We have used the source code discussed in [20]. More details in §2.9 [21].

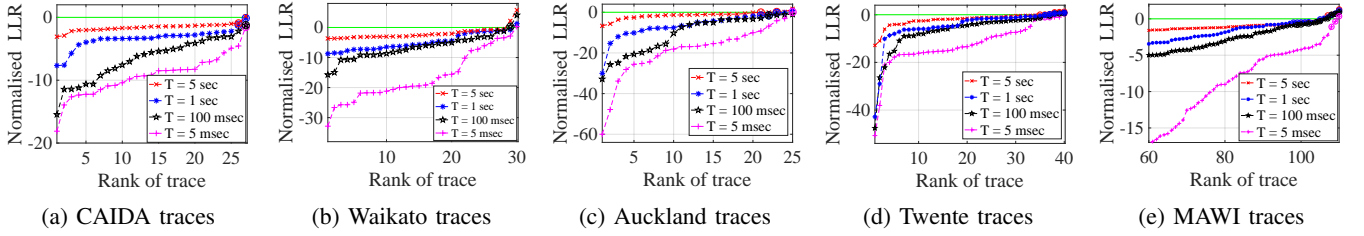


Fig. 2: Normalised Log-Likelihood Ratio ( $\mathcal{R}$ ) test results for all studied traces when running *fit.distCompare*(power-law, log-normal). Aggregation timescales are 5 sec, 1 sec, 100 msec and 5 msec. Circled points in the plots are the ones with  $p$ -value greater than 0.1, i.e. test is inconclusive with respect to assessing which of the reference and candidate distributions is a better fit to the traffic data.

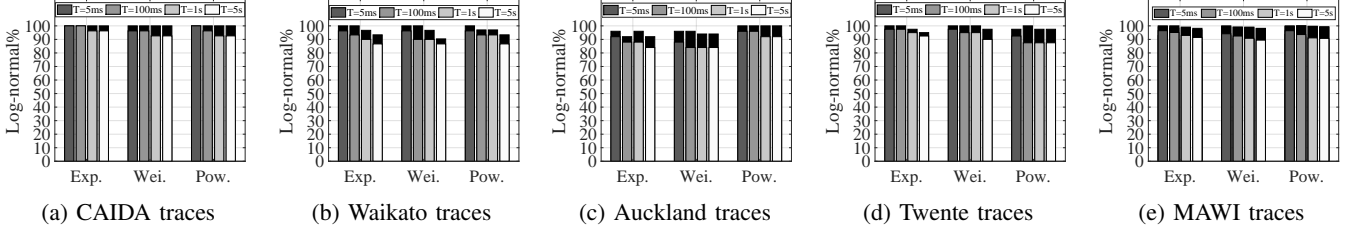


Fig. 3: The percentage of traces for which the log-normal distribution is a better fit compared to the alternative distribution exponential, Weibull or power-law at aggregation timescales: 5 msec, 100 msec, 1 sec and 5 sec. Black areas represent the inconclusive results (i.e.  $p$ -value  $> 0.1$ ).

for the data, and negative if the log-normal distribution is a better fit for the data. The further  $\mathcal{R}$  value is from zero, the better the fit is for one distribution over the other. A  $p$ -value less than 0.1 means that the value of  $\mathcal{R}$  is a reliable indicator of which model (log-normal or alternative, depending on the sign of  $\mathcal{R}$ ) is the better fit to the data. In contrast, a  $p$ -value greater than 0.1 means that there is nothing to be concluded from the likelihood ratio test.

#### A. Fitting the log-normal distribution to Internet traffic data

Figure 1 shows the results of the LLR test for all 232 traces when comparing the log-normal to the exponential, Weibull and power-law distributions. For this test we have aggregated traffic at a timescale  $T = 100$  msec. The points marked with a circle are the ones with  $p > 0.1$ . It is clear that the log-normal distribution is the best fit for the studied traces; i.e.  $\mathcal{R} < 0$  and  $p < 0.1$ <sup>8</sup>. The log-normal distribution is not the best fit for 1 out of 27 CAIDA traces, 3 out of 30 Waikato traces, 3 out of 25 Auckland traces, 5 out of 40 Twente traces and 9 out of 110 MAWI traces i.e. 21 out of 232 are anomalous traces. Most of these traces have inconclusive results (17 out of 21) while for few of them an alternative distribution is a better fit (4 out of 21). We examined these traces in more detail and discuss them in Section III-B.

Identifying the log-normal distribution as the best fit for the vast majority of traffic traces at  $T = 100$  msec is very encouraging. This specific traffic aggregation timescale has been commonly studied in the literature [22], [23].

<sup>8</sup>For clarity, in Figures 1(e) and 2(e) we only plot traces 60 – 110. For traces 1 – 59,  $\mathcal{R}$  is less than 0 and the respective  $p$ -value is less than 0.1; i.e. the log-normal distribution is the best fit for the respective trace.

Next we investigate which one of the studied distributions is the best model for a range of aggregation timescales. We run the fitting test again for the following pairs: (power-law, log-normal), (exponential, log-normal) and (Weibull, log-normal) for timescales between 5 msec to 5 sec. Figure 2 shows the results from Equation 1 of the test between (power-law, log-normal). As reflected by the  $\mathcal{R}$  and  $p$ -values, the log-normal distribution is the best fit for the vast majority of captured traces at all examined timescales<sup>9</sup>. Due to lack of space, we do not show detailed results for the other two pairs; i.e. (exponential, log-normal) and (Weibull, log-normal). Instead, we present overall results, including the percentage of inconclusive tests, for all pairs of distributions, traces and timescales, in Figure 3. The results show that the log-normal is the best fit for the vast majority of traces (overall: 95.19% when  $T=5$  msec and 89.13% when  $T=5$  sec), while few tests are inconclusive (black areas in Figure 3, overall: 4.11% when  $T=5$  msec and 7.44% when  $T=5$  sec) and very few are in favour of the alternative distribution (the rest of the percentages, e.g. overall: 0.7% when  $T=5$  msec and 3.43% when  $T=5$  sec). This is a strong result suggesting the generality of our observations and the potential for wide applicability of the log-normal model in practical applications.

We also examined Q-Q plots for a large number of traces<sup>10</sup>. The log-normal distribution appeared to be a better fit than other tested distributions and no deviations from the expected pattern were observed in the body or tail of the distribution.

<sup>9</sup>Note that it is possible that the network traffic may not follow a log-normal distribution at very fine or coarse aggregation granularities.

<sup>10</sup>Due to lack of space, Q-Q plots are not included as we would have to present plots for each trace, separately.

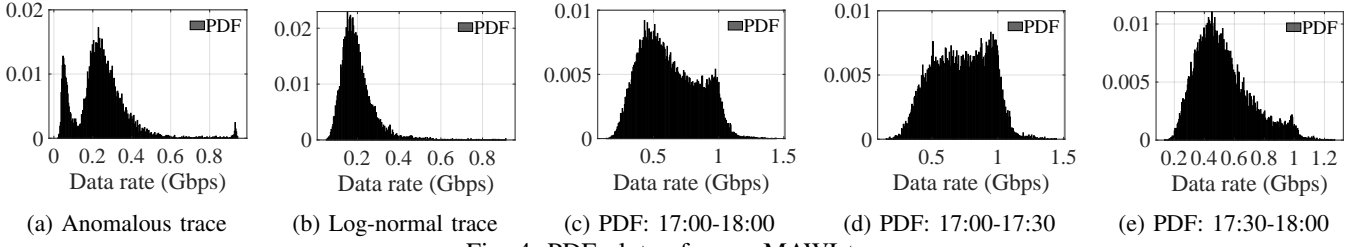


Fig. 4: PDF plots of some MAWI traces.

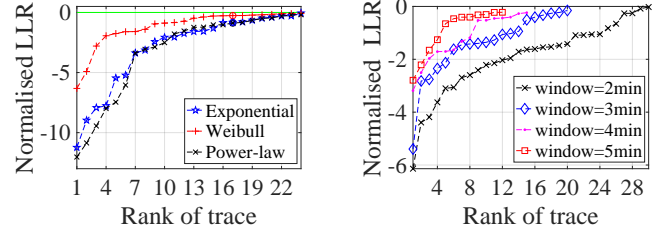
### B. Anomalous traces

As mentioned in Section III-A, there is a small number of traces for which the log-normal distribution is not a good fit (none of the other examined distributions is, either). Figure 4a shows the probability density function (PDF) plot for one of the 9 anomalous MAWI traces. For comparison, Figure 4b shows the PDF for another MAWI trace for which the log-normal distribution is a good fit. It is obvious from Figure 4a that the link was either severely underutilised (see the large spike on the left part of the plot area) or fully utilised (see the smaller spike at the right part of the plot area) for higher data rates. All traces for which the log-normal distribution was not a good fit exhibited similar behaviour and (aggregated) traffic patterns. On the contrary, we did not observe any such behaviour for the majority of traces for which the log-normal distribution was the best fit. A likely explanation for the anomalous traces is that those traces contain either periods of over-capacity (traffic is at 100% of link capacity) or periods where the link is broken (no traffic).

### C. Fitting the log-normal distribution to subtraces in the 24-hour long trace

We need to establish whether we can reliably say that the log-normal distribution is a good fit for any sample length of data, not only the 15-minute long traces (this is discussed in details in Section IV). We apply Clauset test [10] (discussed in Section III) on longer and shorter traces as follows.

Firstly, we apply this test on each subtrace (1-hour long) of the 24-hour long MAWI trace at timescale  $T = 100$  msec. Figure 5a shows the results of the LLR test on 24 subtraces when applying Equation 1. These results complement our results on the 15-minute long traces (see Figure 1) by showing that the log-normal distribution is the best fit for these subtraces. Similar results are seen for other timescales between  $T = 5$  msec and  $T = 5$  sec. There is only one trace (trace id 17 in Figure 5a) where neither the log-normal nor any of the other tested distributions provide a good fit; more specifically, Step 2 of the Clauset method yielded inconclusive results for all distributions. This trace was captured at time 17:00-18:00. The PDF of the 1-hour long trace has two peaks (Figure 4c), and this could be fitted using a bimodal distribution (which we leave as future work). However, when dividing this trace into two 30-minute long subtraces (Figures 4d and e), the log-normal distribution was the best fit for each one of these, separately. Looking at the PDF plots, it is clear that for the first subtrace, the network was much busier compared to the



(a) 24 subtraces (b) windows in 1-hour long trace  
Fig. 5: Normalised LLR test results ( $T = 100$  ms) for MAWI traces (a) 24 subtraces (b) windows in 1-hour trace.

second subtrace. We have no definitive explanation for why this might be, but one could speculate that this is the result of crossing the end of working day or because of some partial equipment failure.

Secondly, we apply the Clauset method on small groups from a 1-hour long MAWI trace. We picked data points from this trace using different time windows: 2, 3, 4 and 5 minutes. This means that each group contains 30, 20, 15 and 12 subtraces, respectively. Figure 5b shows the LLR test results on all these groups when using power-law as alternative to log-normal (same results are seen when using other alternative distributions i.e. exponential and Weibull). The results show that the log-normal distribution is the best fit among all tested distributions for all of these small subtrace groups at all tested windows.

In Section IV we show that the majority of the 15-minute and 1-hour long traces in the dataset are stationary at all tested aggregation timescales. This gives confidence that the log-normal distribution can be used for estimation and modelling purposes as we show in Sections V&VI.

### D. Fitting the log-normal and Gaussian distributions using the correlation coefficient test

The linear correlation coefficient test has been widely used to assess the fit of a distribution to empirical data. To reinforce the results of Section III-A, we assess the fit of the log-normal and Gaussian distributions to all studied traces. We use the linear correlation coefficient as defined in [24]:

$$\gamma = \frac{\sum_{i=1}^n (S_{(i)} - \hat{\mu})(x_i - \hat{x})}{\sqrt{\sum_{i=1}^n (S_{(i)} - \hat{\mu})^2 \cdot \sum_{i=1}^n (x_i - \hat{x})^2}} \quad (2)$$

where  $S_{(i)}$  is the observed sample  $i$ ,  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n S_{(i)}$  is the samples' mean value, and  $x_i$  is sample  $i$  from the



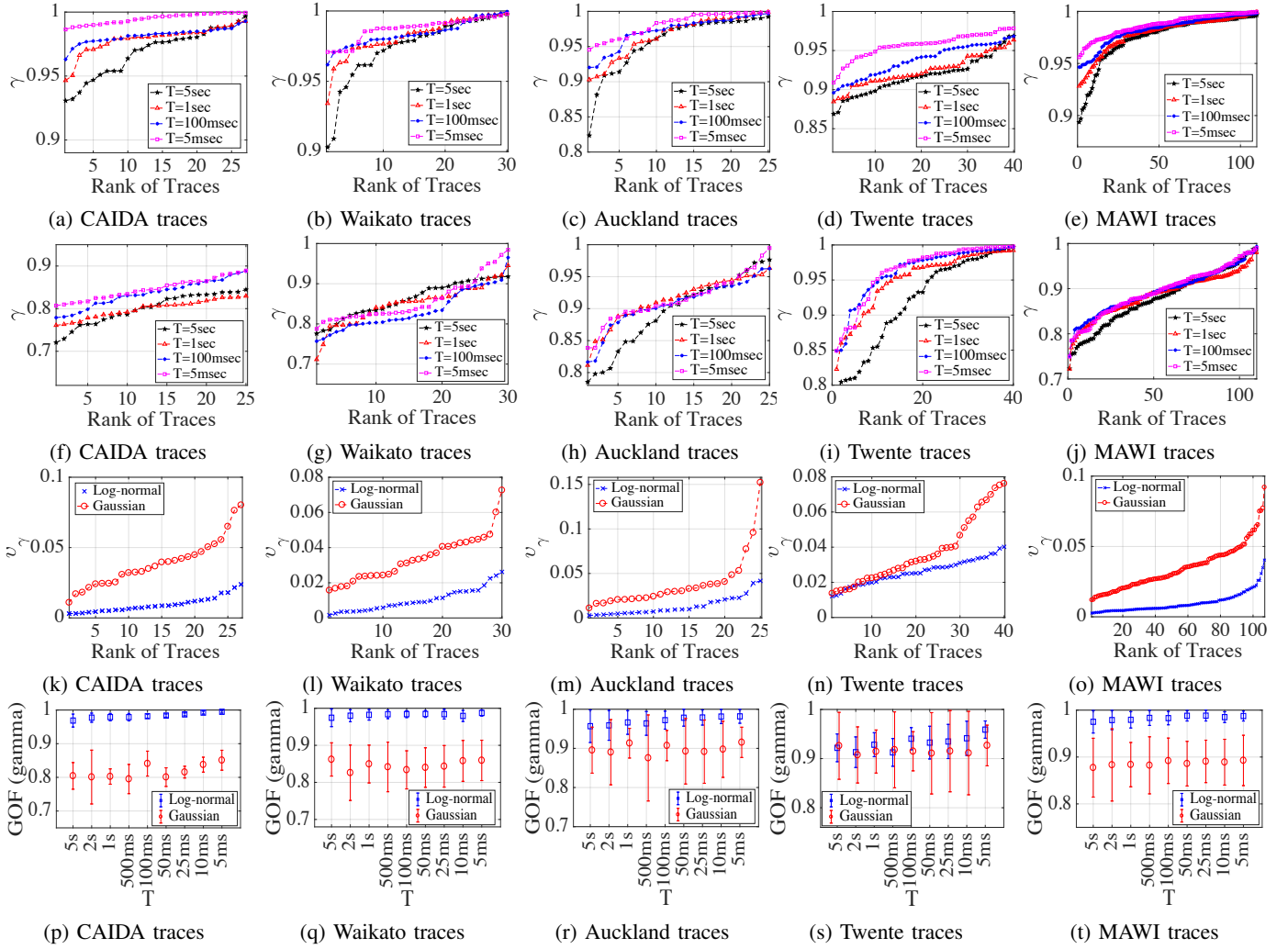


Fig. 6: Correlation coefficient  $\gamma$  test results for all studied traces and different timescales for the log-normal (a-e) and Gaussian (f-j) distributions. The variation  $v_\gamma$  results (k-o). Goodness of fit (GOF) results (p-t).

reference distribution (log-normal in our case), which can be calculated from the inverse cumulative distribution function (CDF) of the reference random variable  $x_i = F^{-1}\left(\frac{i}{n+1}\right)$  and  $\hat{x} = \frac{1}{n} \sum_{i=1}^n x_i$  is the respective mean value. The value of the correlation coefficient can vary between  $-1 \leq \gamma \leq 1$ , with a 1, 0 and  $-1$  indicating perfect correlation, no correlation and perfect anti-correlation, respectively. Strong goodness-of-fit (GOF) is assumed to exist when the value of  $\gamma$  is greater than 0.95 [22].

We measure the linear correlation coefficient for all datasets at four different aggregation timescales (ranging from 5 msec to 5 sec) and plot the results in Figures 6(a-e) for the log-normal distribution and Figures 6(f-j) for the Gaussian distribution. Traces are ordered by the value of  $\gamma$  for the given timescale. It can be clearly seen that  $\gamma > 0.95$  for most traces when employing the test for the log-normal distribution, but this is not the case for the Gaussian distribution.  $\gamma$  is larger for smaller aggregation timescales indicating that the log-normal distribution is an even better fit as the aggregation

gets finer. For very small values of  $T$ , i.e. lower than 1 msec, data samples exhibit binary behaviour, where either a packet is transmitted or not during each examined time frame [23]. We have examined  $\gamma$  for very short (and large) aggregation timescales, and can confirm the absence of a model describing the data (for brevity, we have omitted the relevant figures).

Next, we calculate  $v_\gamma$  (the variation of  $\gamma$ ) for each dataset.  $v_\gamma$  gives an indication of the stability of  $\gamma$  for each dataset, for all timescales tested. This metric is defined as:

$$v_\gamma = \sqrt{\text{var}(\gamma_{T_1}, \gamma_{T_2}, \gamma_{T_3}, \gamma_{T_4})} \quad (3)$$

where  $T_1 = 5\text{ sec}$ ,  $T_2 = 1\text{ sec}$ ,  $T_3 = 100\text{ msec}$  and  $T_4 = 5\text{ msec}$ . Figures 6(k-o) show the results for each dataset with the traces ranked by  $v_\gamma$ . For log-normal model,  $v_\gamma$  is very small (below 0.045) for all traces, therefore we can conclude that  $\gamma$  is almost constant for all studied aggregation timescales. While  $v_\gamma$  is higher for the Gaussian model. Furthermore, the error bars in Figures 6(p-t) represent the standard deviation of the correlation coefficient at different timescales (see x-axis).

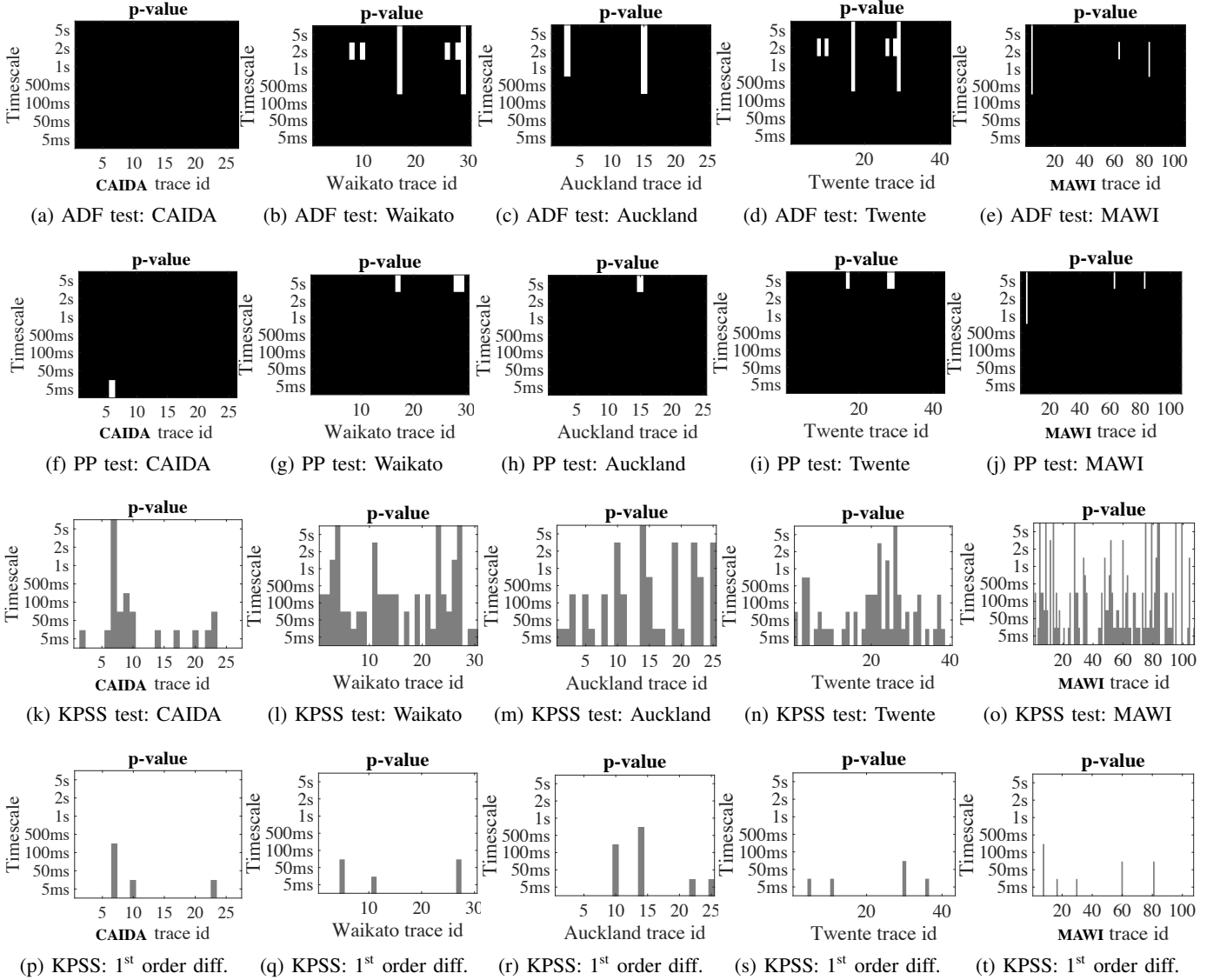


Fig. 7: Stationarity tests' results of the 15-minute long traces. Black: stationary, grey: non-stationary, white: inconclusive. In ADF (a-e) and PP (f-j) tests, the black areas represent  $p\text{-value} \leq 0.05$  (stationary results), while the white areas represent  $p\text{-value} > 0.05$  (inconclusive results). In KPSS test (k-o) and KPSS first-order differencing results (p-t), the grey areas represent  $p\text{-value} \leq 0.05$  (non-stationary results), while the white areas represent  $p\text{-value} > 0.05$  (inconclusive result) (see Table I).

This again shows that for the log-normal model  $\gamma$  is larger than 0.95 (at different  $T$  values) for most CAIDA and MAWI traces, while it is larger than 0.9 for all other datasets. This is not the case with the Gaussian model, where most  $\gamma$  values are less than 0.9.

Overall, the correlation coefficient test reinforces the results extracted in Section III-A, providing strong evidence that the log-normal distribution is the best fit for all studied traces. The superior performance of our model can also be seen from comparison of our results for correlation coefficient with those in [25] where the Gaussian model was used.

#### IV. STATIONARITY TESTING

In the previous section we showed that the log-normal distribution is the best fit among all studied distributions

for the vast majority of traces, for all studied aggregation timescales. In this section, we investigate whether these results are applicable on any traffic sample and we need to find what is the aggregation timescale and time period levels (i.e. length of captured trace) at which these results remain valid. This can be done by applying time series analysis on the traffic by using stationarity tests. From a probabilistic point of view, stationarity means that the distribution remains unchanged when shifted in time. Stationarity plays an important role in time series analysis [12]–[14]. It is important to mention that the study of Internet traffic as a time series depends on two factors. Firstly, a time period over which to study the traffic. Obviously, over long time periods (hours and days) the data is not stationary as it is subject to daily and



TABLE I: ADF, PP and KPSS tests

	ADF and PP tests	KPSS test
null hypothesis (H0)	unit root is present	series is stationary
alternative hypothesis (H1)	series is stationary	unit root is present
$p\text{-value} > 0.05$	H0 is not rejected: result is inconclusive	H0 is not rejected: result is inconclusive
$p\text{-value} \leq 0.05$	H0 is rejected: series is stationary	H0 is rejected: series is non-stationary (due to a unit root)

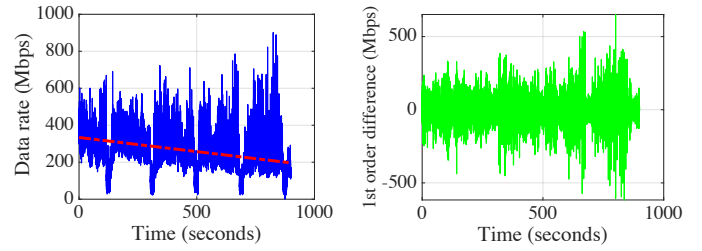
weekly variations related to human activity. Secondly we use a timescale that is used to aggregate the traffic over a specific time period. If the aggregation timescale is very small then the traffic volume will really be a product of exactly how many packets are classified as arriving within that period leading to very noisy measurements. If the timescale is longer, our measured time period will contain very few samples and the statistics calculated will lack power to reject hypotheses, producing instead inconclusive results simply because they have insufficient data. Hence we want to establish whether we can reliably say that a sample of 15-minute or 1-hour long of the captured data is typically stationary in the data set and at which aggregation timescale. Our stationarity test results show that over a 15-minute and 1-hour periods the data is stationary when aggregated at timescales of 0.5 sec to 5 sec.

We examine traffic stationarity using the traces discussed in Section II at different timescales using three tests commonly used for stationarity testing, namely the Augmented Dickey-Fuller (ADF) [26], Phillips-Perron (PP) [27] and Kwiatkowski-Phillips-Schmidt-Shin (KPSS) [28] tests. In the ADF and PP tests, the *null* hypothesis is that a *unit root* is present and the alternate hypothesis is *stationarity*. In the KPSS test, the *null* hypothesis is that the time series is *stationary* and the alternate hypothesis is that a *unit root* is present. Table I summarises the hypotheses of the three tests and outlines the outcome of each test according to the  $p$ -value.

#### A. Stationarity tests of 15-minute long traces

We begin by conducting the stationarity tests on the 232 15-minute long traces (described in Section II). The numbers of data points used in the stationarity tests for each dataset are 180000, 18000, 9000, 1800, 900, 450 and 180 for aggregation timescales of 5 msec, 50 msec, 100 msec, 500 msec, 1 sec, 2 sec and 5 sec, respectively. Figure 7 shows the results for the ADF, PP and KPSS tests at all studied aggregation timescales.

According to the ADF and PP tests (Figures 7(a-j)), the majority of time series are stationary for all aggregation timescales; the  $p$ -value is less than 0.05, therefore the *null* hypothesis is rejected and there is enough evidence to support the alternative hypothesis. These traces are shown as black areas in the figures. There are a few traces for which the  $p$ -value is greater than 0.05 at some aggregation timescales. These are illustrated as white areas in the figures. For these traces the null hypothesis cannot be rejected. These are the *anomalous traces* discussed in Section III-B. Below, we employ the KPSS test to provide evidence that these series are non-stationary; i.e. to show that for the studied traffic traces, where the log-normal was not a good fit for a specific trace, the underlying



(a) Trace with trends

(b) Differenced trace

Fig. 8: First-order differencing of a MAWI trace with trends.

time series was not stationary. As shown in Figures 4a&c, said traces appear with a bi-modal distribution.

In KPSS test results (Figures 7(k-o)), we fail to reject the null hypothesis for most traces, as the  $p$ -value is larger than 0.05. These traces with inconclusive results are shown as white areas in the figures. For some traces (commonly at small aggregation timescales), the null hypothesis is rejected, and therefore there is evidence that the series is non-stationary. These are shown as grey areas in the figures.

**De-trending by differencing in KPSS.** The KPSS test is known to be vulnerable to mistaking trends for non-stationarity i.e. the KPSS test is sensitive to trends [29]. This problem appears for small aggregation timescales because fluctuations appear, and these are mistaken as trends within the time series. If KPSS indicates non-stationarity and ADF indicates stationarity, then the series is difference stationary [29] (that is the case with the grey areas results in Figures 7(k-o)). In order to further explore this observation, we re-ran the analysis by first de-trending the series. De-trending is carried out by using differencing which can be used to remove the series' dependence on time, including structures like trends and seasonality [30]. The differenced series is checked for stationarity as not all non-stationary time series are difference stationary. Figure 8(b) shows the first-order difference of the trace that is shown in Figure 8(a). This gives a  $p$ -value of 0.1, i.e. we fail to reject the stationarity null hypothesis. Figures 7(p-t) show that we always fail to reject the stationarity null hypothesis for most traces (at the tested timescales) in our datasets when applying the KPSS test on the first-order difference sequences (i.e. the results are inconclusive). The results shown in Figure 7(k-t) are consistent with the conclusion that the data is stationary at larger aggregation timescales: 0.5–5 sec and first-difference stationary at smaller aggregation timescales: 5 – 100 msec.

#### B. Stationarity tests of an hour long samples within a 24-hour trace

In this section we consider the hour-long samples within a 24-hour MAWI trace (described in Section II). This 24-hour long trace is used to see if the assumption of stationarity holds for periods longer than 15 minutes. We conclude that these are also stationary.

Figures 9(a-b) show the data rate plots as PDF and as a time series, respectively. It is obvious that the 24-hour long series

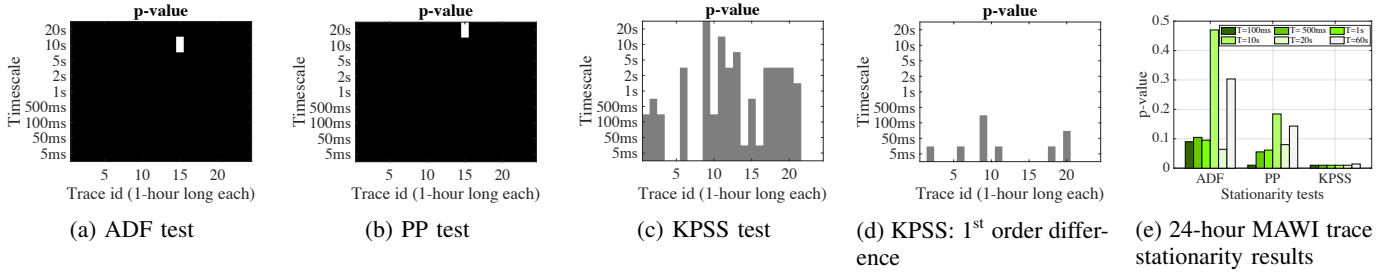
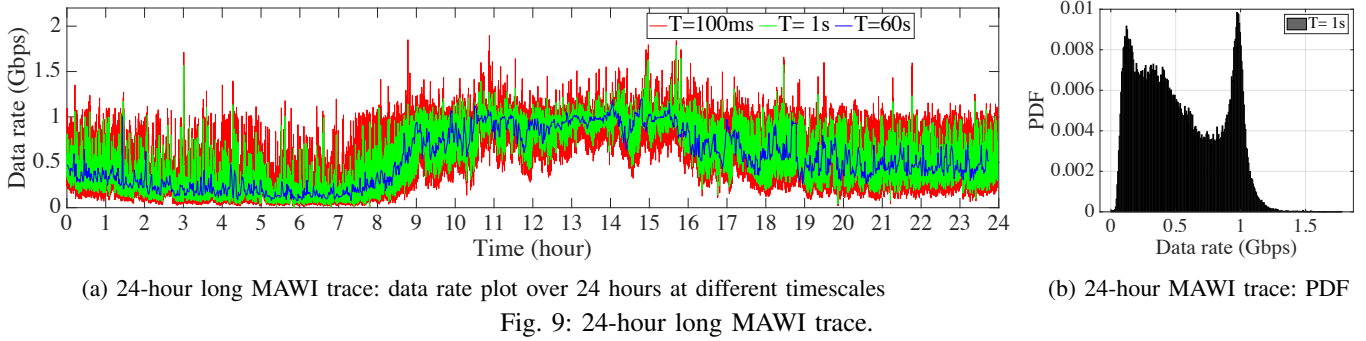


Fig. 10: (a-d) Stationarity tests' results of 24 1-hour long subtraces from the 24-hour long MAWI trace. Black: stationary, grey: non-stationary, white: inconclusive. In ADF and PP tests, the black areas represent  $p\text{-value} \leq 0.05$  (stationary results), while the white areas represent  $p\text{-value} > 0.05$  (inconclusive results). In KPSS test, the grey areas represent  $p\text{-value} \leq 0.05$  (non-stationary results), while the white areas represent  $p\text{-value} > 0.05$  (inconclusive result) (see Table I).

is not stationary; for example, the average data rate in this series between 12:00 am to 05:00 am is 0.252 Gbps, while the average in the time period between 09:00 am to 17:00 pm is 0.875 Gbps.

We run the stationarity tests for this traffic trace and for different aggregation timescales at different sampling times. We start by applying ADF, PP and KPSS tests on 1-hour long groups (subtraces) in this trace (the time at each group or subtrace starts at the beginning of each hour). The stationarity tests results are shown in Figures 10(a-d). It is clear from the first two tests' results that the majority of the 1-hour long groups are stationary (black areas in the figure) as their null hypothesis is rejected. The KPSS test shows many subtraces (white areas) where we failed to reject the null hypothesis i.e. inconclusive results. We ran the KPSS test on the first-order difference series (white areas in Figure 10d) and the number of contradictory results was greatly reduced as the KPSS test failed to reject the null-hypothesis of trend stationarity.

It is worth mentioning that these results might slightly change for some groups if we use 1-hour long groups that do not start at the beginning of each hour (e.g. when using a 1-hour long group that starts at 08:30 am, as a jump in the captured data rate will appear at the second half of this group causing it to be non-stationary). Figure 10e shows the stationarity tests results of the 24-hour long MAWI trace. As expected, and based on the stationarity tests results, this 24-hour long trace is non-stationary.

In the next two sections (V&VI) we present the impact of traffic distribution on two sample traffic engineering problems: link dimensioning and traffic billing. We do not intend our

examples to be fully worked systems for practical deployment. We wish to demonstrate using motivational examples that the improved predictions made possible by these models could in the future have practical utility.

## V. BANDWIDTH PROVISIONING

It has been previously suggested that network link provisioning could be based on fitted traffic models instead of relying on straightforward empirical rules [25]. In this way, over- or under-provisioning can be mitigated or eliminated even in the presence of strong traffic fluctuations. Such approaches rely on having a statistical model that accurately describes the network traffic. This is therefore an excellent area for applying our findings on fitting the log-normal distribution to Internet traffic data. In the literature, the following inequality (the authors call it the “link transparency formula”) has been used for bandwidth provisioning [23]:

$$P(A(T) \geq CT) \leq \varepsilon. \quad (4)$$

In words, this inequality states that the probability that the captured traffic  $A(T)$  over a specific aggregation timescale  $T$  is larger than the link capacity has to be smaller than the value of a performance criterion  $\varepsilon$ . The value of  $\varepsilon$  is chosen carefully by the network provider in order to meet a specific SLA [25]. Likewise, the value of the aggregation time  $T$  should be sufficiently small so that the fluctuations in the traffic can be modelled as well, taking into account the buffering capabilities of network switching devices<sup>11</sup>.

<sup>11</sup>Large traffic fluctuations at very short aggregation timescales are smoothed by the presence of buffers at network routers and switches.

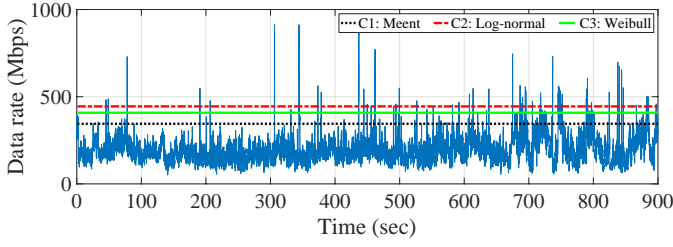


Fig. 11: Data rate of a MAWI trace ( $T = 100$  msec and  $\varepsilon = 0.01$ ). The horizontal lines represent the calculated link capacity based on different models.

We compare bandwidth provisioning using Meent's approximation formula [25] (assuming Gaussian) and using a log-normal traffic model.

#### A. Bandwidth provisioning using Meent's formula

To find the minimum required link capacity, Meent et al. [25] proposed a bandwidth provisioning approach that is based on the assumption that the traffic follows a Gaussian distribution. Meent's dimensioning formula is defined as follows [25]:

$$C1 = \mu + \frac{1}{T} \sqrt{-2 \log(\varepsilon) \cdot v(T)} \quad (5)$$

where  $\mu$  is the average value of the traffic,  $v(T)$  is the variance at timescale  $T$  and  $\varepsilon$  is the performance criterion. The link capacity is obtained by adding a safety margin value

$$\text{Safety margin} = \sqrt{-2 \log(\varepsilon)} \cdot \sqrt{\frac{v(T)}{T^2}}$$

to the average of the captured traffic (see Equation 5). This safety margin value depends on  $\varepsilon$  and the ratio  $\sqrt{v(T)/T^2}$ . As the value of  $\varepsilon$  decreases the safety margin increases. For example, when the value of  $\varepsilon$  decreases from  $10^{-2}$  to  $10^{-4}$ , then value of the safety margin increases by 40%. This is different from conventional link dimensioning methods, where the safety margin is fixed to be 30% above the average of the presented traffic [25], [31], [32]. Traffic tails are represented using the Chernoff bound, as follows:

$$P(A(T) \geq CT) \leq e^{-SCT} E[e^{SA(T)}]. \quad (6)$$

Here  $E[e^{SA(T)}]$  is the moment generation function (MGF) of the captured traffic  $A(T)$ .

#### B. Bandwidth provisioning based on the log-normal model

Here we investigate whether we could achieve more reliable bandwidth provisioning by adopting the log-normal traffic model. We calculate the mean and variance from the captured trace and generate the respective log-normal model. Then, we use the CDF function ( $F$ ) to solve the link transparency formula shown in Equation 4. Hence,  $F$  is defined as  $F(C) = P(A(T)/T < C)$ , which can be solved to find  $C$ , as follows:

$$C2 = F^{-1}(1 - \varepsilon). \quad (7)$$

#### C. Comparison of bandwidth provisioning approaches

In this section, we compare the bandwidth provisioning approaches described above. The performance indicator is the empirical value of the performance criterion, which is denoted by  $\hat{\varepsilon}$  and defined as follows:

$$\hat{\varepsilon} = \frac{\#\{A_i | A_i \geq CT\}}{n}, i \in 1 \dots n. \quad (8)$$

In words, this empirical value is the percentage of all the data samples of the captured traffic which are measured larger than the estimated link capacity. Ideally,  $\hat{\varepsilon}$  would be equal to the target value of the performance criterion  $\varepsilon$ . The difference between  $\hat{\varepsilon}$  and  $\varepsilon$  is due to the fact that the chosen traffic model is not accurately describing the real network traffic. A simple example of the described comparison approach is illustrated in Figure 11, in which we plot the captured data rate for a MAWI trace ( $T = 100$  msec)<sup>12</sup>. The calculated capacity values from each approach when the target  $\varepsilon$  is 0.01 are  $C1 = 344.8$  Mbps and  $C2 = 444.3$  Mbps (represented by the horizontal lines in Figure 11). The empirical value can be calculated by using Equation 8, which gives  $\hat{\varepsilon}_1 = 0.042$  and  $\hat{\varepsilon}_2 = 0.012$ . Obviously, with the first approach the network operator would not be able to meet the target  $\varepsilon = 0.01$ , while with the second approach the empirical value is close to the target.

We next compare results of bandwidth provisioning calculations based on the (a) Meent's formula, (b) Weibull model and (c) proposed log-normal model. Figure 12(a-d) shows the average of the empirical value ( $\text{avg}(\hat{\varepsilon})$ ) for all traces in each dataset at  $T = 0.1$  sec,  $T = 0.5$  sec and  $T = 1$  sec. The value of  $T$  is chosen to be sufficiently small so that the fluctuations in the traffic can be modelled as well. Each model is tested for four different values of the performance criterion:  $\varepsilon = 0.5$ ,  $\varepsilon = 0.1$ ,  $\varepsilon = 0.05$  and  $\varepsilon = 0.01$ . In Figure 12(a-d) we clearly see that the log-normal model is able to satisfy the required performance criterion  $\varepsilon$  at different aggregation time-scales for all datasets. In contrast, Meent's formula failed to allocate sufficient bandwidth, which results in missing the target performance criterion  $\varepsilon$  for all datasets and target performance values, as depicted in Figure 12(i-l) (see horizontal red line). The Weibull distribution performs better comparing to Meent's formula, but bandwidth provisioning using the log-normal model is far superior, as can be seen from Figures 12(a-d) and 12(e-h).

We apply the same link dimensioning tests as discussed above on 24 subtraces (each one being 1-hour long) from the 24-hour long MAWI trace. Figure 13 shows the  $\text{avg}(\hat{\varepsilon})$  for all subtraces at different timescale values. As shown in the figure, the log-normal model performs the best compared to the other two in estimating bandwidth allocation, with respect to the target performance criterion.

Over long time periods (hours and days) the data is not stationary as it is subject to daily and weekly variations related to human activity. A single distribution cannot sensibly capture

<sup>12</sup>Note that in all subsequent figures we have also included results for a Weibull model to get insights about bandwidth provisioning using a heavy-tailed distribution.

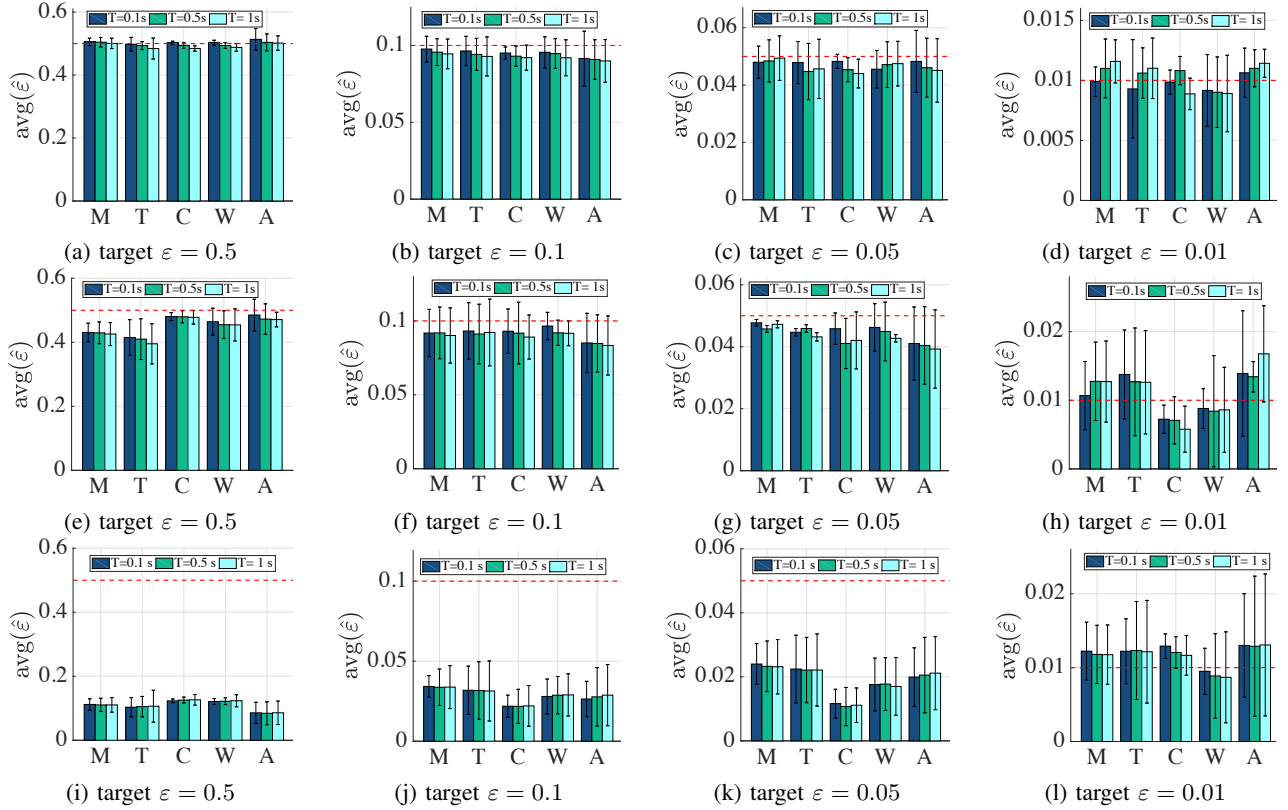


Fig. 12: Link dimensioning based on (a-d) log-normal model, (e-h) Weibull model and (i-l) Meent's formula:  $\text{avg}(\hat{\varepsilon})$  for different datasets (M: MAWI, T: Twente, C: CAIDA, W: Waikato, A: Auckland), aggregation timescales (100 msec, 500 msec and 1 s), and target values of  $\varepsilon$  (0.5, 0.1, 0.05 and 0.01). Error bars represent  $\text{stderr} |\varepsilon - \hat{\varepsilon}|$ .

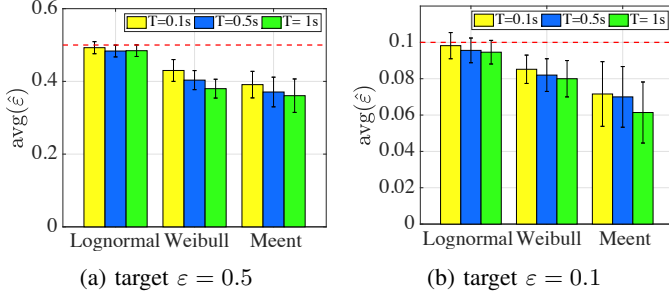


Fig. 13: Link dimensioning based on log-normal, Weibull and Meent for 24 subtraces from 24-hour long MAWI trace.

the behaviour of traffic over one day. In Figure 5a we showed that for all 1-hour long traces the log-normal distribution is the best fit compared to all tested alternative distributions, therefore, a sensible procedure for the operator could be to fit individual log-normal distributions to smaller periods of the day. The traffic would then be modelled as a series of log-normal distributions where the mean and variance change between times of day. A practical indicative procedure that an operator can follow to use the log-normal model for bandwidth provisioning and provide the respective analysis can be as in the following example. We divide the Mawi 24-hour long trace into 96 15-minute long subtraces starting from time 00:00 to 23:59. Then, we apply the bandwidth provisioning mechanism

on these 96 samples by measuring the empirical value  $\hat{\varepsilon}$  for two target values:  $\varepsilon = 0.5$  and  $\varepsilon = 0.1$  using different models (log-normal, Weibull and Meent). Figure 14 shows that the log-normal model is able to achieve the target performance values much better than Weibull and Meent models (NRMSE = 0.0396 and 0.0052 for targets  $\varepsilon = 0.1$  and  $\varepsilon = 0.5$ , respectively, in log-normal results). These results show that the log-normal model can accurately predict the proportion of time a link will exceed a given capacity. As an example, the operator could choose to provision bandwidth based on the fitted model for the peak time of the day. It is worth pointing out that we do not intend our example to be fully worked systems for practical deployment (see future work in Section VIII).

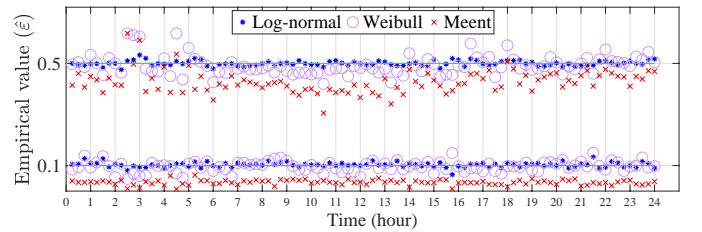


Fig. 14: Link dimensioning for 96 15-minute subtraces from the 24-hour Mawi trace for two target  $\varepsilon$  values.



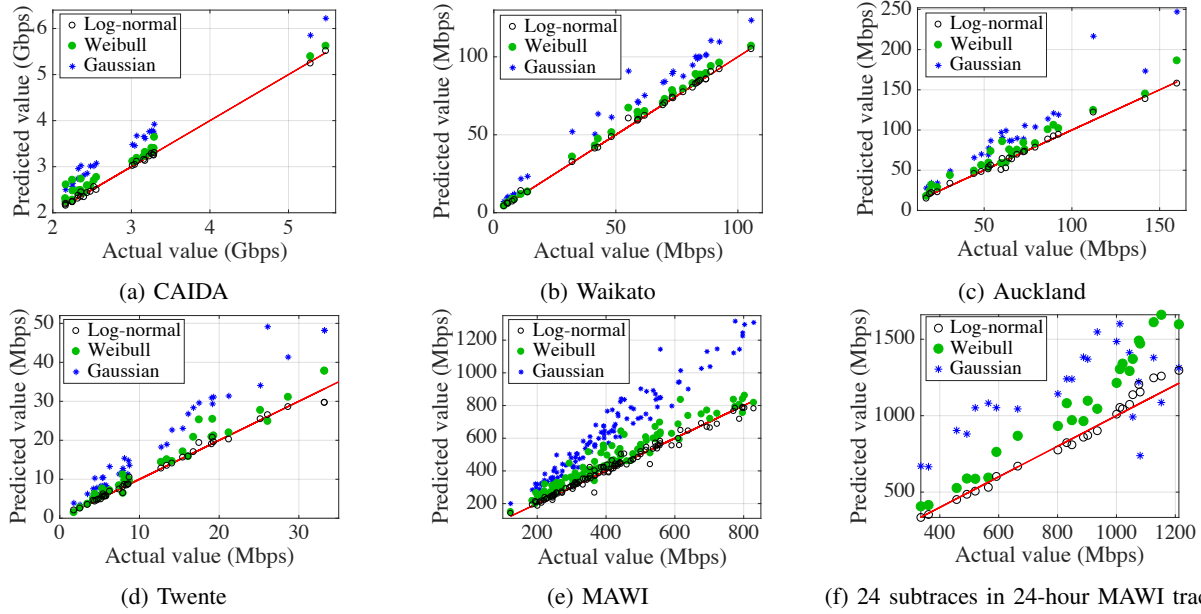


Fig. 15: 95th percentile values (actual vs predicted rates) based on log-normal, Weibull and Gaussian models. An ideal model would result in points in the plot area that fall exactly on the red line. The results in (a-e) are for the 15-minute long traces in each dataset. The results in (f) are for the 24 subtraces in the 24-hour long MAWI trace.

## VI. 95TH PERCENTILE PRICING SCHEME BASED ON LOG-NORMAL MODEL

Traffic billing is typically based on the 95th percentile method [33]. Traffic volume is measured at border network devices (typically aggregated at time intervals of 5 minutes) and bills are calculated according to the 95th-percentile of the distribution of measured volumes; i.e. network operators calculate bills by disregarding occasional traffic spikes. Forecasting future bills, which is important for ISPs and clients, can be done using a model of the traffic calculated through previously sampled traffic. In this section, we apply our findings on Internet traffic modelling in predicting the cost of traffic according to the 95th percentile method.

For each network trace we calculate the actual 95th percentile of the traffic volume. The majority of the studied traffic traces were 15-minute long but operators typically use measurements traffic volumes for much longer periods, therefore we scale down the calculation of the 95th percentile by dividing each trace (900 seconds) into 90 groups (10 seconds length each). In reality, of course, the 95th percentile method would use traffic from different times of day with different means which would need to be modelled as separate log normal distributions with separate means and variances. This is not possible with the 15-minute long samples that form the focus of this paper. However, it remains a level playing field test for which of the distributions best captures the real underlying data.

We calculate the 95th percentile for the observed traffic. We then fit a Gaussian, Weibull and log-normal distribution to each trace (for  $T = 100$  msec) and calculate the 95th percentile of the fitted distribution. We plot the actual 95th percentile

against the three predictions in Figure 15 with a red reference line to show where perfect predictions would be located. It is clear that the log-normal model provides much more accurate predictions of the 95th percentile than the Gaussian model. As with the bandwidth dimensioning case discussed in Section V, the Weibull is better than the Gaussian model but worse than the proposed log-normal model.

We employ the normalised root mean squared error (NRMSE) as a goodness of fit to the results in Figure 15. NRMSE measures the differences between values predicted by a hypothetical model and the actual values. In other words, it measures the quality of the fit between the actual data and the predicted model. Table II shows the NRMSE for all datasets and the three considered models. It is clear that the lowest NRMSE value is for the log-normal model, which is the best model compared to the Gaussian and Weibull ones.

A simple procedure of how an operator might use the log-normal model in predicting the 95th percentile would be to fit individual log-normal distributions to smaller periods of the day. The traffic is then being modelled as a series of log-normal distributions. We apply the same 95th percentile test on each 1-hour long subtrace from the 24-hour long MAWI trace. For each 1-hour long subtrace, we calculated the 95th percentile by dividing each subtrace (1 hour) into 60 groups (1-minute long each). Figure 15f shows the results for the 24 subtraces. The log-normal model is significantly more accurate in predicting the 95th percentile compared to the Weibull and Gaussian models. This example is not intended to be a fully deployable system. It aims at highlighting the benefits of our findings. Several months of data would be necessary to investigate the applicability of this model (see future work in Section VIII).

TABLE II: Goodness of fit (GOF) using Normalised Root Mean Squared Error (NRMSE)

Model/Dataset	CAIDA	Waikato	Auckland	Twente	MAWI
Log-normal	0.0399	0.0401	0.1058	0.0979	0.1528
Weibull	0.2410	0.1148	0.2984	0.2123	0.4145
Gaussian	0.5544	0.4193	0.6866	0.5741	0.9828

## VII. RELATED WORK

Reliable traffic modelling is important for network planning, deployment and management; e.g. for traffic billing and network dimensioning. Historically, network traffic has been widely assumed to follow a Gaussian distribution. In [5], [7], the authors studied network traces and verified that the Gaussianity assumption was valid (according to simple goodness-of-fit (GOF) tests they used) at two different timescales. In [34], the authors studied traffic traces during busy hours over a relatively long period of time and also found that the Gaussian distribution is a good fit for the captured traffic. Schmidt et al. [8] found that the degree of Gaussianity is affected by short and intensive activities of single network hosts that create sudden traffic bursts. All the above mentioned works agreed on the Gaussian or ‘fairly Gaussian’ traffic at different levels of aggregations in terms of timescale and number of users. The authors in [24], [35] examined the levels of aggregation required to observe Gaussianity in the modelled traffic, and concluded that this can be disturbed by traffic bursts. The work in [9], [36] reinforces the argument above, by showing existence of large traffic spikes at short timescales which result in high values in the tail. Compared to existing literature, our findings are based on a modern, principled statistical methodology, and traffic traces that are spatially and temporally diverse. We have tested several hypothesised distributions and not just Gaussianity.

An early work drawing attention to the presence of heavy tails in Internet file sizes (not traffic) is that of Crovella and Bestavros [2]. Deciding whether Internet flows could be heavy-tailed became important as this implies significant departures from Gaussianity. The authors in [37] provided robust evidence for the presence of various kinds of scaling, and in particular, heavy-tailed sources and long-range dependence in a large dataset of traffic spanning a duration of 14 years.

When modelling network traffic, many authors did not perform any tests for stationarity, assuming that their traces are realisations of a weakly stationary stochastic process [38], [39]. Other authors consider non-stationary behaviour of traffic observations [14], [40]. The authors in [41], [42] show that traffic patterns have almost deterministic daily variations resulting in clear non-stationary behaviour on a day timescale. The authors in [40] demonstrated that multiplexed traffic on a high-speed link may have non-stationary behaviour and discussed possible causes of non-stationarity of traffic observations. They argue that this could be due to time-varying number of aggregated sources, routing changes or specific aggregation of a constant number of stationary sources. A common approach in traffic modelling is to choose sufficiently small blocks of observations such that observations in separate

blocks are expected to be at least weakly stationary [13], [22], [43]. For example, when testing for applicability of the Gaussian model to traffic modelling authors in [24] neglected a part of their trace claiming that it may introduce undesirable non-stationary behaviour. Authors in [14] assumed that 5-minute blocks of their traffic observations are sufficient to ensure intra-block stationarity.

Understanding the traffic characteristics and how these evolve is crucial for ISPs for network planning and link dimensioning. Operators typically over-provision their networks. A common approach to do so is to calculate the average bandwidth utilisation [6] and add a safety margin. As a rule of thumb, this margin is defined as a percentage of the calculated bandwidth utilisation [31]. Meent et al. [25] proposed a new bandwidth provisioning formula, which calculates the minimum bandwidth that guarantees the required performance, according to an underlying SLA. This approach relies on the statistical parameters of the captured traffic and a performance parameter. The underlying fundamental assumption for this to work is that the traffic the network operator sees follows a Gaussian distribution. Same approach has been used in [23].

The 95th percentile method is used widely for network traffic billing. Dimitropoulos et al. [33] have found that the computed 95th percentile is significantly affected by traffic aggregation parameters. However, in their approach they do not assume any underlying model of the traffic; instead, they base their study on specific captured traces. Stanojevic et al. [4] proposed the use of Shapley value for computing the contribution of each flow to the 95th percentile price of interconnect links. Works [44]–[47] propose calculating the 95th percentile using experimental approaches. Xu et al. [48] assume that network traffic follows a Gaussian distribution “through reasonable aggregation” and propose a cost-efficient data centre selection approach based on the 95th percentile.

## VIII. CONCLUSION

The distribution of traffic on Internet links is an important fundamental problem that has received relatively little attention. We use a well-known, state-of-the-art statistical framework to investigate the problem using a large corpus of traces. The traces cover several network settings including home user access links, tier 1 backbone links and campus to Internet links. The traces are from times from 2002 to 2020 and are from a number of different countries. We investigated the distribution of the amount of traffic observed on a link in a given (small) aggregation timescale which we varied from 5ms to 5s. The hypotheses compared were that the traffic volume was heavy-tailed, that the traffic was log-normal and that the traffic was normal (Gaussian). The vast majority of traces fitted the log-normal assumption best and this remained true for all timescales tried. Where no distribution tested was a good fit this could be attributed either to the link being saturated (at full capacity) for a large part of the observation or exhibiting signs of link-failure (no or very low traffic for part of the observation).



We tested the data for the hypothesis of stationarity. Over long periods (hours and days) the data is not stationary as it is subject to daily and weekly behaviour related to human activity. Over a fifteen minute or one hour period our tests show that the data is stationary when aggregated at timescales of 500ms to 5s and is first-difference stationary when aggregated at smaller time-scales from 5ms to 100ms.

We investigate the impact of the distribution on two sample traffic engineering problems. Firstly, we looked at predicting the proportion of time a link will exceed a given capacity. This could be useful for provisioning links or for predicting when SLA violation is likely to occur. Secondly, we looked at predicting the 95th percentile transit bill that ISP might be given. For both of these problems the log-normal distribution gave a more accurate result than a heavy-tailed distribution or a Gaussian distribution. We conclude that the log-normal distribution is a good (best) fit for traffic volume on normally functioning internet links in a variety of settings and over a variety of timescales, and further argue that this assumption can make a large difference to statistically predicted outcomes for applied network engineering problems.

**Limitations and future work.** It is important to point out that the presented procedures on bandwidth provisioning and 95th-percentile pricing are not meant to be fully worked systems for practical deployment. This would require testing our model with larger data sets that last for days, weeks or even months. Instead, we intended to motivate the need for deriving good data models for Internet traffic volumes which would be crucial in developing real-world systems. As part of our future work, we will investigate the possibility of developing a bimodal distribution that can fit the anomalous traces and explore more candidate distributions for fitting Internet traffic volume data. Also, studying several months of data would be necessary to investigate the variability between hours of the day, days of the week and so on.

## REFERENCES

- [1] P. Pruthi *et al.*, “Heavy-tailed on/off source behavior and self-similar traffic,” in *Proc. of ICC*, 1995.
- [2] M. E. Crovella *et al.*, “Self-similarity in World Wide Web traffic: evidence and possible causes,” in *IEEE/ACM ToN*, 1997.
- [3] P. Loiseau *et al.*, “Investigating Heavy-Tailed Distributions on a Large-Scale Experimental Facility,” in *IEEE/ACM ToN*, 2010.
- [4] R. Stanojevic *et al.*, “On Economic Heavy Hitters: Shapley Value Analysis of 95Th-percentile Pricing,” in *Proc. of IMC*, 2010.
- [5] R. Meent *et al.*, “Gaussian traffic everywhere?” in *Proc. of ICC*, 2006.
- [6] R. d. O. Schmidt *et al.*, “Measurement-based network link dimensioning,” in *Proc. of IFIP Networking*, 2015.
- [7] R. d. O. Schmidt, R. Sadre *et al.*, “Gaussian traffic revisited,” in *Proc. of IFIP Networking*, 2013.
- [8] R. d. O. Schmidt, R. Sadre, N. Melnikov *et al.*, “Linking network usage patterns to traffic Gaussianity fit,” in *Proc. of IFIP Networking*, 2014.
- [9] X. Yang, “Designing traffic profiles for bursty Internet traffic,” in *Proc. of GLOBECOM*, 2002.
- [10] A. Clauset *et al.*, “Power-law distributions in empirical data,” *arXiv:0706.1062v2*, 2009.
- [11] M. Alasmar *et al.*, “On the Distribution of Traffic Volumes in the Internet and its Implications,” in *Proc. of INFOCOM*, 2019.
- [12] G. Lauks *et al.*, “Testing the Null Hypothesis of Stationarity of Internet Traffic,” in *Elektronika Ir Elektrotechnika*, 2011.
- [13] D. Moltchanov, “Modeling local stationary behavior of Internet traffic,” in *Journal of Communications Software and Systems*, 2008.
- [14] J. Cao *et al.*, “On the Nonstationarity of Internet Traffic,” in *Proc. of SIGMETRICS*, 2001.
- [15] “The CAIDA UCSD Anonymized Internet Traces,” 2016. [Online]. Available: [http://www.caida.org/data/passive/passive\\_dataset.xml](http://www.caida.org/data/passive/passive_dataset.xml)
- [16] “Mawi Archive,” 2018. [Online]. Available: <http://mawi.wide.ad.jp/>
- [17] R. Barbosa *et al.*, “Simpleweb/University of Twente Traffic Traces Data Repository,” <http://eprints.eemcs.utwente.nl/17829/>, Tech. Rep., 2010.
- [18] “WITS: Waikato Internet Traffic Storage,” 2013. [Online]. Available: <https://wand.net.nz/wits/waikato/8/>
- [19] “WITS: Auckland X,” 2009. [Online]. Available: <https://wand.net.nz/wits/auck/10/>
- [20] J. Alstott *et al.*, “powerlaw: a Python package for analysis of heavy-tailed distributions,” in *arXiv:1305.0215*, 2014.
- [21] M. Alasmar, “Understanding the characteristics of Internet traffic and designing an efficient RaptorQ-based data transport protocol for modern data centres,” in *Ph.D. thesis, University of Sussex*, <https://sro.sussex.ac.uk/id/eprint/89371/>, 2019.
- [22] M. Mandjes and R. van de Meent, “Resource Dimensioning Through Buffer Sampling,” in *IEEE/ACM Transactions on Networking*, 2009.
- [23] R. d. O. Schmidt *et al.*, “Impact of Packet Sampling on Link Dimensioning,” in *Transactions on Network and Service Management*, 2015.
- [24] J. Kilpi *et al.*, “Testing the Gaussian Approximation of Aggregate Traffic,” in *Proc. of SIGCOMM*, 2002.
- [25] A. Pras *et al.*, “Dimensioning network links: a new look at equivalent bandwidth,” in *IEEE Network*, 2009.
- [26] D. Dickey and W. Fuller, “Distribution of the Estimators for Autoregressive Time Series With a Unit Root,” in *JSTOR*, 1979.
- [27] P. C. B. Phillips and P. Perron, “Testing for a Unit Root in Time Series Regression,” in *Journal of the American Statistical Association*, 1988.
- [28] D. Kwiatkowski *et al.*, “Testing the null hypothesis of stationarity against the alternative of a unit root,” in *Journal of Econometrics*, 1992.
- [29] S. Seabold *et al.*, “statsmodels: Econometric and statistical modeling with python,” in *9th Python in Science Conference*, 2010.
- [30] R. Hyndman *et al.*, “Forecasting: Principles and Practice,” in *OTexts: [Online] https://otexts.com/fpp2/*, 2018.
- [31] “Best Practices in Core Network Capacity Planning,” online, accessed December 2020. [Online]. Available: [https://www.cisco.com/c/en/us/products/collateral/routers/wan-automation-engine/white\\_paper\\_c11-728551.html](https://www.cisco.com/c/en/us/products/collateral/routers/wan-automation-engine/white_paper_c11-728551.html)
- [32] M. Alasmar *et al.*, “Network link dimensioning based on statistical analysis and modeling of real internet traffic,” *arXiv:1710.00420*, 2017.
- [33] X. Dimitropoulos *et al.*, “On the 95-Percentile Billing Method,” in *Proc. of PAM*, 2009.
- [34] García-Dorado *et al.*, “Characterization of the busy-hour traffic of IP networks based on their intrinsic features,” in *Computer Networks*, 2011.
- [35] A. B. Downey, “Evidence for Long-tailed Distributions in the Internet,” in *Proc. of SIGCOMM Workshop on Internet Measurement*, 2001.
- [36] H. Abrahamsson *et al.*, “Traffic characteristics on 1Gbit/s access aggregation links,” in *Proc. of ICC*, 2017.
- [37] R. Fontugne and *et al.*, “Scaling in internet traffic: A 14 year and 3 day longitudinal study, with multiscale analyses and random projections,” *IEEE/ACM Transactions on Networking*, 2017.
- [38] W. E. Leland *et al.*, “On the self-similar nature of Ethernet traffic,” in *IEEE/ACM ToN*, 1994.
- [39] W. Willinger *et al.*, “Self-similarity through high-variability: statistical analysis of Ethernet LAN traffic,” in *IEEE/ACM ToN*, 1997.
- [40] T. Karagiannis *et al.*, “A nonstationary Poisson view of Internet traffic,” in *Proc. of INFOCOM*, 2004.
- [41] K. Thompson *et al.*, “Wide-area Internet traffic patterns and characteristics,” in *IEEE Network*, 1997.
- [42] V. Paxson *et al.*, “Wide area traffic: the failure of Poisson modeling,” in *IEEE/ACM Transactions on Networking*, 1995.
- [43] M. Jainik *et al.*, “Measurement and modelling of the temporal dependence in packet loss,” in *Proc. of INFOCOM*, 1999.
- [44] L. Golubchik and *et al.*, “To send or not to send: Reducing the cost of data transmission,” in *Proc. of INFOCOM*, 2013.
- [45] N. Laoutaris *et al.*, “Inter-datacenter Bulk Transfers with Netstitcher,” in *Proc. of SIGCOMM*, 2011.
- [46] I. Castro *et al.*, “Using Tuangou to Reduce IP Transit Costs,” in *IEEE/ACM Transactions on Networking*, 2014.
- [47] H. Xu *et al.*, “Joint request mapping and response routing for geo-distributed cloud services,” in *Proc. of INFOCOM*, 2013.
- [48] —, “Cost efficient datacenter selection for cloud services,” in *Proc. of ICC*, 2012.